



# Optimization of Gene Expression by Natural Selection

## Citation

Bedford, Trevor and Daniel L. Hartl. 2009. Optimization of gene expression by natural selection. *Proceedings of the National Academy of Sciences* 106(4): 1133-1138.

## Published Version

<http://dx.doi.org/10.1073/pnas.0812009106>

## Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:3630584>

## Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

# Optimization of gene expression by natural selection

Trevor Bedford<sup>1</sup> and Daniel L. Hartl<sup>2</sup>

Department of Organismic and Evolutionary Biology, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138

Contributed by Daniel L. Hartl, November 25, 2008 (sent for review October 20, 2008)

It is generally assumed that stabilizing selection promoting a phenotypic optimum acts to shape variation in quantitative traits across individuals and species. Although gene expression represents an intensively studied molecular phenotype, the extent to which stabilizing selection limits divergence in gene expression remains contentious. In this study, we present a theoretical framework for the study of stabilizing and directional selection using data from between-species divergence of continuous traits. This framework, based upon Brownian motion, is analytically tractable and can be used in maximum-likelihood or Bayesian parameter estimation. We apply this model to gene-expression levels in 7 species of *Drosophila*, and find that gene-expression divergence is substantially curtailed by stabilizing selection. However, we estimate the selective effect,  $s$ , of gene-expression change to be very small, approximately equal to  $Ns$  for a change of one standard deviation, where  $N$  is the effective population size. These findings highlight the power of natural selection to shape phenotype, even when the fitness effects of mutations are in the nearly neutral range.

evolution | nearly neutral | Ornstein-Uhlenbeck | phenotypic optima

**A**bundant evidence indicates that natural selection is remarkably powerful in shaping nucleotide sequences (1, 2). Many tests of natural selection rely on a comparison between nonsynonymous sites, in which mutations affect protein sequence, and synonymous sites, in which mutations do not. Synonymous sites serve as a proxy for neutral sites, enabling the effects of selection to be distinguished from background mutational and demographic patterns. Although changes in gene expression are hypothesized to play a major role in adaptation (3, 4), changes in expression cannot be so easily partitioned into neutral and selected categories. Thus, methods derived to analyze selection in coding sequences cannot be readily applied to gene-expression data. In part because of this ambiguity, general forces acting on gene-expression divergence have remained unclear. At this point, there exists considerable debate over the relative importance of selection and random drift in shaping gene-expression levels (5–8).

The benefits of optimal gene regulation seem in many ways obvious. In the simple case of metabolic enzymes, under-expression may slow metabolic flux, while over-expression may expose the cell to additional toxic misfolded proteins (9). At the morphological level, gene regulation can be tightly coupled to phenotype (10, 11). Genetic mutations whose effects cascade into morphological differences are expected to have especially large fitness impacts, and as such will be heavily influenced by natural selection. A straightforward example of selection on gene-expression level can be seen in ribosomal proteins, which contrary to the neutral prediction are found to be highly expressed across a variety of organisms (12).

In this article, we present a model of gene-expression divergence that explicitly distinguishes between the forces of random genetic drift and natural selection. This work is based upon prior models of phenotypic trait evolution (13, 14). Our population genetic model is fundamentally similar to the Brownian motion model used to describe the random movements of physical particles (15). In both cases, the system is impacted by numerous tiny perturbations, in Brownian motion caused by collision but

in the evolutionary context caused by mutations that are fixed in an evolving population. Owing to the central limit theorem, the resulting state of the system can be accurately described as a normally distributed random variable. In the simplest case, the probability of fixation of a random mutation is assumed to be independent of the current state of the system, and thus movement is not favored in one direction over the other. This scenario corresponds to selective neutrality. However, a slightly more complex model, described by the Ornstein-Uhlenbeck (OU) process, assumes that perturbations are more likely to shift the system toward some optimal value than away from it (16). This model does well to capture the essence of natural selection; mutations that produce a phenotype closer to some optimum are favored over those that produce a phenotype farther away.

Here, we analyze gene-expression levels across 7 species of *Drosophila* using the framework provided by the OU model. In the analysis, we compare expression divergence between species with estimates of time since their divergence based on sequence data. The pattern at which divergence in gene-expression levels accumulates over time does much to reveal the underlying forces of selection and drift. Using only species-level data, we find that stabilizing selection plays a major role in limiting divergence of gene-expression level. We also quantify the degree of selection and drift for specific genes, which illuminates the relationship between changes in gene sequence and changes in gene expression. Finally, we reconstruct the fitness landscape of gene-expression level, and find that although natural selection is pervasive in shaping gene expression, the individual fitness effects of changes in gene expression are rather weak.

## Modeling Expression Divergence

**Analogy to Brownian Motion.** Here we apply models of Brownian motion to describe the variance in gene-expression level between orthologous genes as a function of the time separating these orthologs (13, 14). Brownian motion, also known the Wiener process, represents one of simplest continuous-time, continuous-state stochastic processes. In a Brownian motion, the degree of stochastic change away from the current state is independent of both state and time. The increment that a Brownian motion makes over a time interval of length 1 is normally distributed with mean 0 and variance  $\sigma^2$ . The “volatility” parameter  $\sigma$  completely describes the Brownian motion and determines the rate at which a trait’s value diffuses away from its current state. In an evolutionary context,  $\sigma$  describes that rate of “phenotypic drift” experienced by a gene. Our use of the term drift differs from the classic usage, wherein drift refers to a systematic trend in the evolution of a Brownian motion. Genes in which expres-

Author contributions: T.B. and D.L.H. designed research; T.B. performed research; T.B. contributed new reagents/analytic tools; T.B. and D.L.H. analyzed data; and T.B. and D.L.H. wrote the paper.

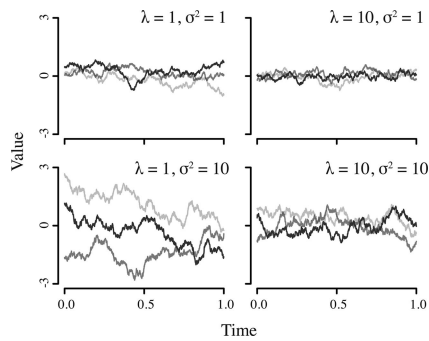
The authors declare no conflict of interest.

<sup>1</sup>To whom correspondence may be addressed at: Department of Ecology and Evolutionary Biology, University of Michigan, 2041 Kraus Natural Science Building, 830 North University, Ann Arbor, MI 48109. E-mail: bedfordt@umich.edu.

<sup>2</sup>To whom correspondence may be addressed. E-mail: dhartl@oeb.harvard.edu.

This article contains supporting information online at [www.pnas.org/cgi/content/full/0812009106/DCSupplemental](http://www.pnas.org/cgi/content/full/0812009106/DCSupplemental).

© 2009 by The National Academy of Sciences of the USA



**Fig. 1.** Realizations of the OU process. Three individual realizations are shown for each of four different parameter values. The drift parameter  $\sigma$  determines the degree of mutational pressure randomly impacting the trait value, while  $\lambda$  determines the pull of selection toward some optimal trait value (in this case 0). In each realization, the starting value was sampled from the equilibrium distribution.

sion has a larger mutational target size (17) are expected to show larger values of  $\sigma$ . The probability density function of a Brownian motion is:

$$f(x|x_0, \sigma, t) \sim \mathcal{N}(x_0, t \sigma^2)$$

where  $x_0$  is equal to the state of the process at time 0. Thus, Brownian motion predicts that the extent of variance in gene-expression increases in proportion to time. This scenario corresponds to selective neutrality, as the model assumes that change in expression is independent of current expression level.

Selection favoring an optimal level of gene expression can be incorporated using a simple extension to the Brownian motion model (13, 14, 18). This addition results in an OU or mean-reverting process (16). If Brownian motion is thought of as a particle that is subject to random perturbations from its surroundings, then an OU process can be thought of as adding an elastic spring to this particle, attaching it at some fixed point. As random perturbations push the particle farther away from this fixed point, the strength of elastic return increases proportionally. Thus, in addition to the stochastic force of drift, an OU process includes the deterministic force of selection pulling the trait toward some optimal value. The instantaneous motion of an OU process is described by:

$$dx = (\mu - x) \lambda dt + \mathcal{N}(0, \sigma dt)$$

where  $\mu$  represents the optimal trait value,  $\lambda$  is proportional to the strength of selection, and  $\sigma$  is proportional to the strength of drift. Solving this yields the density function of an OU process:

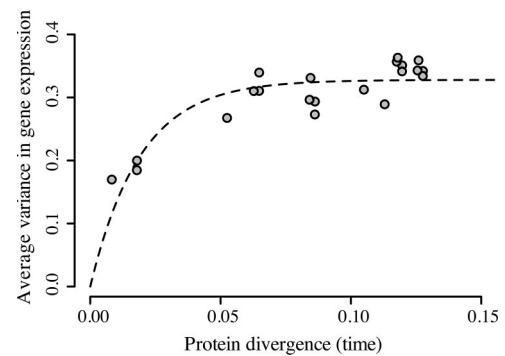
$$f(x|x_0, \mu, \lambda, \sigma, t) \sim \mathcal{N}\left(x_0 e^{-\lambda t} + \mu (1 - e^{-\lambda t}), \frac{\sigma^2}{2\lambda} (1 - e^{-2\lambda t})\right)$$

Here we see that variance does not increase in proportion to time, and instead saturates at a stable equilibrium:

$$\lim_{t \rightarrow \infty} f(x|x_0, \mu, \lambda, \sigma, t) \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{2\lambda}\right)$$

The temporal character of the OU model for various values of  $\lambda$  and  $\sigma$  is shown in Fig. 1.

**Inferring Fitness Landscapes.** We convert the OU parameters  $\lambda$  and  $\sigma$  into population-genetic estimates of the strength of selection through comparison of the ratio of the instantaneous rates of positive change and negative change in the OU model to the ratio of fixation rates of selectively advantageous and disadvantageous



**Fig. 2.** Average pairwise variance in expression level for *Drosophila* species. Each point represents the average variance between a species pair. This variance initially increases with time, but eventually saturates. In the absence of stabilizing selection, pairwise variance is expected to saturate at 1. Non-linear regression fit of pairwise variance vs. time for the OU model is represented as a dashed line ( $\lambda = 26.14$ ;  $\sigma = 4.14$ ).

mutations. We find that the ratio of instantaneous rates of change for the OU model is:

$$\frac{r_{x \rightarrow y}}{r_{y \rightarrow x}} = \frac{\lambda}{\sigma^2} (x - y)(x + y - 2\mu)$$

Following Kimura (19), we find the ratio of fixation rates between mutants of  $+Ns$  and  $-Ns$  effect to be:

$$\frac{r_+}{r_-} = \left( \frac{2Ns}{1 - e^{-2Ns}} \right) / \left( \frac{-2Ns}{1 - e^{-2Ns}} \right) = \frac{e^{2Ns} - 1}{1 - e^{-2Ns}} = e^{2Ns}$$

Here, the equation is simplified by multiplying numerator and denominator by  $e^{2Ns}$ . Thus, the rate difference between positive and negative change in the OU model can be used to derive an  $Ns$  value by setting these two equations equal to each other and solving for  $Ns$ :

$$Ns_{x \rightarrow y} = \log \sqrt{\frac{r_{x \rightarrow y}}{r_{y \rightarrow x}}} = \frac{\lambda}{2\sigma^2} (x - y)(x + y - 2\mu)$$

If we measure relative to the optimum (i.e., fitness at optimum = 1), then this expression reduces to  $Ns(z) = 1 - z^2\lambda/2\sigma^2 = 1 - z^2/4v$ , where  $z$  represents the distance to the optimum in terms of standard deviations, and  $v$  represents expected equilibrium variance. Thus, the curvature of the fitness landscape is inversely proportional to the level of equilibrium variance observed. As such, we will refer to equilibrium variance as measuring the degree of selective constraint that the expression level of a gene experiences. It is this measure of selective constraint rather than the  $\lambda$  parameter that should be used in comparing selection across genes or across species, as the observed value of  $\lambda$  depends upon both selective constraint and mutational input.

## Results

One key finding is that the accumulation of variance in gene-expression level between 7 species of *Drosophila* is not proportional to the amount of time separating each species (Fig. 2). This result immediately suggests that continuous neutral evolution of gene expression is unlikely. Instead, we find that expression divergence between orthologous genes saturates rapidly in evolutionary time. This general pattern was previously hypothesized to exist by Whitehead and Crawford (20). Species pairs of *Drosophila* do not show a significant increase in expression divergence beyond that present between *D. melanogaster* and *D. ananassae*. Saturation of gene-expression divergence is expected if expression levels are under stabilizing selection.





levels. These findings highlight the “overwhelming odds against the less fit” (23) and the power of natural selection to shape phenotypic variation.

The extent of stabilizing selection on gene-expression divergence has been a contentious topic. Khaitovich *et al.* (5), using a similar approach to the present study, find that pairwise divergence in expression level increases in proportion to time across primates. The discrepancy between these results and our own may come from multiple sources. Khaitovich *et al.* examine chimpanzee, orangutan, and macaque expression levels using probes designed for human genes. In this case, sequence differences among species will mimic expression divergence (7), and so apparent expression divergence will continue to increase with time, even when the underlying expression divergence has saturated. Additionally, Khaitovich *et al.* define expression divergence as squared mean difference between species-specific expression levels. This statistic (unlike our measure of average variance, mean of one half of squared differences) is biased by an amount proportional to sampling variance. Phylogenetically distant comparisons had a smaller sample size than close comparisons and so were biased toward large estimates of expression divergence (7). Another study of primate-expression divergence using species-specific probes found that, in the majority of cases, a constant level gene expression across the phylogeny could not be rejected (24). Although this result is consistent with stabilizing selection, a low rate of neutral divergence will have the same effect. Other studies using various methodologies have suggested that stabilizing selection acts upon expression divergence (25–28). However, identifying stabilizing selection in these studies has relied on information in addition to species-specific expression levels. The OU model provides a simple framework for investigating stabilizing selection that requires only expression data from orthologous genes. The OU model allows the degree of stabilizing selection to be compared not only between genes but also between organisms.

**Mutational Input and Genetic Drift.** Random genetic drift eventually results in the conversion of standing genetic variation into fixed differences. We find that empirical estimates of the rate of phenotypic drift in expression level are remarkably consistent with expected rates of random genetic drift, given levels of standing variation and effective population size. Phenotypic drift results in  $\sigma^2 = 17.14$  units of variance in the time it takes to accumulate 1.0 aa substitutions per site. This is equivalent to  $8.68 \times 10^{-10}$  units of expression variance per generation (see *Methods*). Lande (13) gives the expected variance per generation because of random genetic drift as  $h^2\pi^2/N$ , where  $h^2$  is the heritability of the trait,  $\pi^2$  is the level of variance across individuals within a population, and  $N$  is the effective population size. Assuming  $h^2 = 0.5$ ,  $\pi^2 = 0.0726$  (based upon empirical comparisons between two strains of *D. simulans*), and  $N = 9.05 \times 10^6$  [determined from synonymous genetic diversity in *D. simulans* (29) and inferred *Drosophila* mutation rate (30)], we arrive at an expectation of  $4.02 \times 10^{-9}$  units of variance per generation. The reasonably close correspondence between the empirical estimate and the theoretical prediction suggests that the OU model does well to describe the underlying evolutionary process.

However, mutation-accumulation experiments have suggested much larger values of mutational variance in gene-expression level, or  $\approx 2.4 \times 10^{-5}$  units of variance per generation (31). In this study, a relatively small number of individual mutations resulted in widespread changes in gene-expression level. This discrepancy can be reconciled by assuming that mutations of large effect would be purged by natural selection before reaching appreciable frequency and, hence, do not end up contributing to standing genetic variation. This phenomenon is another aspect of selective constraint. Our calculated rate of phenotypic drift of

$\approx 10^{-9}$  represents the population-level turnover of standing variation into fixed differences, and not the input of variation because of new mutations.

**Model Assumptions.** Our analysis has made several simplifying assumptions, including constant gene-expression optima, symmetrical mutation rates, and strong-selection/weak-mutation dynamics. If the optimum itself is subject to stochastic variation, then our analysis will underestimate the true strength of stabilizing selection. This is because movement of the optimum and subsequent tracking by natural selection will appear similar to weak selection poorly tracking a constant optimum. However, strong selection tracking a shifting optimum will result in decreased levels of standing variation compared to levels expected under a constant optimum. We find levels of within-population variation that are highly compatible with the observed rate of drift, suggesting that shifting optima have not had a major influence on our results.

We find that asymmetrical mutation should not significantly impact our results. We simulated evolution on the fitness landscape shown in Fig. 3 under a strong-selection/weak-mutation model, where the rate of mutation to lower expression was twice the rate of mutation to higher expression. We found that asymmetrical mutation had no discernable effect on equilibrium variance (Fig. S3), suggesting our estimates are robust to the presence of mutational asymmetry. Additionally, the results of Lande (13) suggest that our model is robust to the assumption of strong-selection/weak-mutation dynamics.

Throughout our analysis, we have assumed that species-specific normalization (see *Methods*) had little effect on our estimates of OU parameters. To assess the impact of this assumption, we performed simulations wherein expression levels of 10,000 genes were evolved according to the OU model and subsequently normalized in a species-specific fashion (Fig. S4). We find that normalization results in overestimation of the degree of selective constraint, suggesting that our conclusion of nearly neutral evolution is conservative.

## Conclusions

It is well known that purifying selection constrains the rate of sequence change. Often, the reduction in evolutionary rate estimated using  $d_N/d_S$  is taken as a measurement of the degree of selective constraint. We find that selection, rather than simply decreasing the overall rate of expression divergence, instead curtails expression divergence in a nonlinear fashion. Thus, measurement of selective constraint on the evolution of continuous traits requires comparison of multiple orthologous trait values to be successful, but fortunately does not require a neutral proxy in the way of sequence evolution.

The OU framework presented here may be substantially extended to model further intricacies of gene-expression evolution. For example, large-scale fluctuations in  $\lambda$  and  $\sigma$  could be investigated by allowing branch-specific parameter values. We would expect fluctuations of effective population size to significantly impact inferred levels of selection. Additionally, it is possible to identify lineage-specific adaptation for a particular gene by allowing for multiple trait optima across a phylogeny (i.e.,  $\mu$  of *D. melanogaster* may differ from  $\mu$  of other *Drosophila*). Standard methods, such as likelihood-ratio tests, could then be used to assess significance. It would be highly interesting to see whether lineages undergoing adaptive-sequence evolution also show evidence of adaptive gene-expression evolution. We believe that the OU model presented here will prove useful to the future study of gene-expression evolution, and to the study of phenotypic evolution in general.



22. Ohta T (1992) The nearly neutral theory of molecular evolution. *Annu Rev Ecol Syst* 23:263–286.
23. Wallace AR (1892) Note on sexual selection. *Nat Sci* 1:749–750.
24. Gilad Y, Oshlack A, Smyth GK, Speed TP, White KP (2006) Expression profiling in primates reveals a rapid evolution of human transcription factors. *Nature* 440:242–245.
25. Oleksiak MF, Churchill GA, Crawford DL (2002) Variation in gene expression within and among natural populations. *Nat Genet* 32:261–266.
26. Rifkin SA, Kim J, White KP (2003) Evolution of gene expression in the *Drosophila melanogaster* subgroup. *Nat Genet* 33:138–144.
27. Lemos B, Meiklejohn CD, Cáceres M, Hartl DL (2005) Rates of divergence in gene expression profiles of primates, mice, and flies: stabilizing selection and variability among functional categories. *Evolution* 59:126–137.
28. Whitehead A, Crawford DL (2006) Neutral and adaptive variation in gene expression. *Proc Natl Acad Sci USA* 103:5425–5430.
29. Begun DJ, et al. (2007) Population genomics: whole-genome analysis of polymorphism and divergence in *Drosophila simulans*. *PLoS Biol* 5:e310.
30. Tamura K, Subramanian S, Kumar S (2004) Temporal patterns of fruit fly (*Drosophila*) evolution revealed by mutation clocks. *Mol Biol Evol* 21:36–44.
31. Rifkin SA, Houle D, Kim J, White KP (2005) A mutation accumulation assay reveals a broad capacity for rapid evolution of gene expression. *Nature* 438:220–223.
32. *Drosophila* 12 Genomes Consortium (2007) Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* 450:203–218.
33. Yang Z (1997) PAML: a program package for phylogenetic analysis by maximum likelihood. *CABIOS* 13:555–556.
34. Zhang Y, Sturgill D, Parisi M, Kumar S, Oliver B (2007) Constraint and turnover in sex-biased gene expression in the genus *Drosophila*. *Nature* 450:233–237.
35. Oshlack A, Chabot AE, Smyth GK, Gilad Y (2007) Using DNA microarrays to study gene expression in closely related species. *Bioinformatics* 23:1235–1242.